

Universidade Federal de Uberlândia

Bacharelado em sistemas de informação

Faculdade de computação

Organização e recuperação da informação

Atividade prática TP2

Uberlândia – MG

Universidade Federal de Uberlândia

Bacharelado em sistemas de informação

Faculdade de computação

Maycon douglasBatista dos santos11921BSI209

Uberlândia - MG

Relatório

Documento de relatório para o TP2 referentes aos modelos booleanos e $tf*idf$ mostrando os resultados obtidos durante a prática de **organização e recuperação da informação** na disciplina curso **Bacharelado em Sistemas de Informação**.

Professor: Rodrigo Sanches Miani

Sumário

- 1.....
- 2.....
- 3.....
- 4.....
- 5.....
- 6.....
 - 1.....
 - 2.....
- 7.....
 - 1.....
 - 2.....
 - 3.....
- 8.....
- 9.....

Objetivo

O objetivo desta aula é entender revisar o modelo booleano aplicando o conceito de frequência de termos e a partir desse entendimento poder aplicar conceitos mais avançados de recuperação da informação que são usados até hoje.

Exercício 1)

Criar um programa que faça o seguinte:

1) Leia múltiplos arquivos de texto (um diretório, por exemplo)

```
6 # ler todos os documentos de uma pasta com a extensão indicada e
7 # salva na pasta de destino indicada
8 def bagOfWords2 (source = './docs', ext = '.txt' , dest = None):
9
10     bag = []
11     docs = []
12
13     for dir, subPastas, arquivos in os.walk(source):
14
15         for arquivo in arquivos:
16
17             if Regex.findall(ext + '$', arquivo):
18
19                 arg = open(source + '/' + arquivo, 'r')
20
21                 arg = uteis.normaliza(arg.read())
22
23                 if(dest != None):
24                     uteis.saveFile(arg, dest + '/' + arquivo)
25
26                 docs.append(arquivo)
27                 bag.append(arg)
28
29     return bag, docs
```

2) Crie o vocabulário (termos de indexação únicos) dos arquivos de texto removendo qualquer tipo de pontuação e considerando somente letras minúsculas

```
36
37 def normaliza(words):
38
39     bag = []
40
41     palavras = words.replace('\n', ' ').replace(',', '').split(' ')
42
43     for item in palavras:
44         elem = unicode.decode(item.lower())
45
46         if elem not in bag and elem != '':
47             bag.append(elem)
48
49     bag.sort()
50     return bag
51
```

3) Grave o vocabulário em arquivo texto

```

52 # recebe as lista de palavras ja normalizada e grava im arquivo
53 # com extensão .index
54
55 def saveFile (bagOfWords, name):
56
57     arg = open(name + '.index','w')
58
59     for item in bagOfWords:
60         arg.write(item + '\n')
61
62

```

4) Leia o vocabulário em arquivo texto e exiba no console a BoW dos documentos (ausência/presença do termo em um determinado documento)

```

49 # ler todos os docs de uma pasta e salva a bag of words da cada um em uma pasta index
50 if __name__ == '__main__':
51
52     source = input('Entre com o caminho relativo dos arquivos que serão indexados: ')
53
54     ext = input('Qual extensão dos arquivos que serão indexados: ')
55
56     dest = input('Entre com o caminho relátivo de onde serão guardados: ')
57
58     listaBOW, namesFiles = bagOfWords2 ( source , ext, dest)
59
60     name = input('Agora entre com o arquivo com a busca: ')
61
62     keyWords = uteis.normaliza(open(name).read())
63
64     result = matrixBollean(listaBOW, keyWords)
65
66     exhibirBOW(result, namesFiles, keyWords)
67

```

```

maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$ python3 Ex2/ex2.py
entre com nome do arquivo vocabulário: keywords.txt
entre com nome/caminho do diretório: docs
['to', 'do', 'is', 'be', 'or', 'not', 'i', 'am', 'what', 'think', 'therefore', 'da', 'let', 'it']
[0.0, 1.073, 0.0, 0.0, 0.0, 0.0, 2.0, 1.0, 0.0, 2.0, 2.0, 0.0, 0.0, 0.0] -> doc3.txt
[3.0, 0.83, 4.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0] -> doc1.txt
[2.0, 0.0, 0.0, 0.0, 2.0, 2.0, 2.0, 2.0, 2.0, 0.0, 0.0, 0.0, 0.0, 0.0] -> doc2.txt
[0.0, 1.073, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 5.17, 4.0, 4.0] -> doc4.txt
maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$ ls
docs Ex1 Ex2 Ex3 hinos index keywords.txt relatorio.odt
maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$

```

Exercício 2

Criar um programa que receba como parâmetro um arquivo com o vocabulário e um diretório ou caminho com os documentos (arquivos de texto). O programa, então, deverá:

1) Calcular e exibir no console o TF-IDF de cada documento. Testar com a coleção de quatro documentos do livro/sala de aula. Verificar se os valores estão corretos.

```
51 | # calcula a frequencia de uma palavra na lista passasa como primeiro
52 | # parametro (vocabulario indexado) na lista do segundo parametro (doc indexado)
53 | def frequencia (keyWords, lista):
54 |
55 |     result = []
56 |     for key in keyWords:
57 |
58 |         result.append(lista.count(key))
59 |
60 |     return result
61 |
```

```
61 |
62 |
63 | # calcula a frequencia da lista de vocabulario (segundo parametro) na matrix
64 | # que seria o nossa bag (segundo parametro) usando a função descrita acima
65 | def TF (bag, keyWords):
66 |
67 |     result = []
68 |
69 |     for lista in bag:
70 |
71 |         result.append(frequencia(keyWords, lista))
72 |
73 |     return result
74 |
```

```
75 | # recebe um TF e converte para sua variação
76 | def TFvar (tf):
77 |
78 |     result = []
79 |     for doc in tf:
80 |         linha = []
81 |         for item in doc:
82 |
83 |             if(item > 0):
84 |                 aux = 1 + math.log2(item)
85 |             else:
86 |                 aux = 0
87 |
88 |             linha.append(aux)
89 |         result.append(linha)
90 |
91 |     return result
92 |
```



```

93 # calcula o idf para cada termo (coluna) na minha matriz de frequencia
94 def IDF (tf):
95
96     idf = []
97     for i in range(len(tf[0])):
98
99         idf.append(0)
100
101     for i in range(len(tf)):
102         for j in range(len(tf[i])):
103
104             if(tf[i][j] > 0):
105                 idf[j] += 1
106
107     for i in range(len(idf)):
108         if idf[i] > 0:
109             idf[i] = math.log2(len(tf) / idf[i])
110
111     return idf

```

```

113 # pega cada valor da matriz e multiplica pelo seu valor idf
114 def TF_IDF (tf, idf):
115
116     result = []
117     for lin in range(len(tf)):
118
119         linha = []
120         for col in range(len(tf[lin])):
121
122             linha.append(round(float(tf[lin][col]) * float(idf[col]), 3))
123
124         result.append(linha)
125
126     return result
127

```

```

128 if __name__ == '__main__':
129
130     keyWords = open(input("entre com nome do arquivo vocabulário: ")).read()
131
132     keyWords = normaliza(keyWords, ord=False)
133
134     nomeDir = input("entre com nome/caminho do diretório: ")
135
136     bag, nomes = bagOfWords(source=nomeDir)
137
138     tf = TF(bag, keyWords)
139
140     tf = TFvar(tf)
141
142     idf = IDF(tf)
143
144     final = TF_IDF(tf, idf)
145
146     print(keyWords)
147     for i in range(len(final)):
148
149         print(f'{final[i]} -> {nomes[i]}')
150

```

Resultado:

```

maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$ ls
docs Ex1 Ex2 Ex3 hinos index keywords.txt relatorio.odt
maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$ python3 Ex2/ex2.py
entre com nome do arquivo vocabulário: keywords.txt
entre com nome/caminho do diretório: docs
['to', 'do', 'is', 'be', 'or', 'not', 'i', 'am', 'what', 'think', 'therefore', 'da', 'let', 'it']
[0.0, 1.073, 0.0, 0.0, 0.0, 0.0, 2.0, 1.0, 0.0, 2.0, 2.0, 0.0, 0.0, 0.0] -> doc3.txt
[3.0, 0.83, 4.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0] -> doc1.txt
[2.0, 0.0, 0.0, 0.0, 2.0, 2.0, 2.0, 2.0, 2.0, 0.0, 0.0, 0.0, 0.0, 0.0] -> doc2.txt
[0.0, 1.073, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 5.17, 4.0, 4.0] -> doc4.txt
maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$

```

Exercício 3

1) Crie uma coleção de documentos com os hinos dos clubes dos 20 times que irão jogar a Série A do Campeonato Brasileiro de 2023.

```

maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$ ls hinos/
'Hino do Botafogo.txt'      'Hino do Goiás Esporte Clube.txt'      'Hino Oficial do Clube Atlético Bragantino.txt'
'Hino do Corinthians.txt'  'Hino do Grêmio.txt'                  'Hino Oficial do Clube Atlético Mineiro.txt'
'Hino do Cruzeiro.txt'    'Hino do Internacional.txt'            'Hino Oficial do Coritiba Foot-ball Club.txt'
'Hino do Cuiabá Esporte Clube (MT).txt' 'Hino do Palmeiras.txt'                'Hino Oficial do Esporte Clube Bahia.txt'
'Hino do Flamengo.txt'    'Hino do Vasco da Gama.txt'            'Hino Oficial do Santos Futebol Clube.txt'
'Hino do Fluminense.txt'  'Hino Oficial do América Futebol Clube (MG).txt' 'Hino Oficial do São Paulo Futebol Clube.txt'
'Hino do Fortaleza Esporte Clube.txt'  'Hino Oficial do Club Athletico Paranaense.txt'
maycon@maycon-Inspiron-3584:~/5Periodo/ORI/TP2$

```

2) Use o programa criado no Exercício 2 para responder as seguintes perguntas:

Código:

```

4  # recebe um doc (td*idf) valor do menor peso e nomes dos arquivos lidos
5  # considerando que o doc e o nomes tem o mesmo tamanho e a
6  # posição i nos nomes indicada se refere a no doc pelo valor (tf*idf)
7  def menorPeso (doc, valor, nomes):
8
9      lista = []
10     for i in range(len(doc)):
11
12         if valor == doc[i]:
13             lista.append(nomes[i])
14     return lista
15

```

```

17  # apresentação do exe 3
18  if __name__ == '__main__':
19
20     keyWords = open(input("entre com nome do arquivo vocabulário: ")).read()
21
22     nomeDir = input("entre com nome/caminho do diretório: ")
23
24     bag, nomes, voc = uteis.Ex2(keyWords, nomeDir)
25
26     print(f'O tamanho do vocabulário é : {len(voc)}')
27
28     print(f'vocabulário: {voc}')
29     for nome in nomes:
30         print(f'doc {nome} tf*idf é :')
31
32         for doc in bag:
33             print(doc)
34
35             menor = min(doc)
36
37             lista = menorPeso(doc, menor, voc)
38
39             print(f'menor valor {menor} com os termos {lista}')
40

```

a) Qual o tamanho do vocabulário da coleção?

Linha 26

b) Encontre o TF-IDF de cada um dos documentos da coleção.

Linha 24

c) Qual termo(s), em qual documento, possui o maior peso TF-IDF?

Linhas 37, 39

Resultado:

[illegible]